

Using Corpus Statistics in the Modeling of Linguistic Paradigms

Beata Trawiński

SFB 441 – University of Tübingen, Germany

`trawinski@sfs.uni-tuebingen.de`

This paper presents how corpus statistics can be used to verify complex inflectional paradigms. This will be demonstrated using a set of traditionally assumed inflectional paradigms of third person personal pronouns in Polish.

According to traditional approaches, the inventory of third person personal pronouns in Polish comprises entities specified in (1)–(5). Four grammatical categories play a role in constituting the inflectional paradigm of the pronouns, i.e., case (nominative, genitive, dative, accusative, instrumental and locative), number (singular and plural), postprepositionality (yes or no; see forms separated by a slash in (1)–(5)) and accentability (yes or no; indicated by pronoun forms in parentheses versus forms not in parentheses in (1)–(5)). In (1)–(5), only genitive and accusative masculine singular pronouns are specified as being able to appear in the unaccented postprepositional form (cf. the form *-ń* in (1), (2) and (3)). However, the corpus data seem to demonstrate more possibilities (for our investigations, the IPI PAN Corpus of Polish was used, cf. <http://korpus.pl>).

Table 1 presents the distribution of selected unaccented postprepositional pronouns found in the corpus. For each form, the context in which it occurs is specified, i.e., the contraction of that form with a particular preposition, and the total number of times it occurred together with the percentage of the total frequency of all unaccented postprepositional forms is recorded. In addition, the total of all occurrences of each contraction found in the corpus is indicated, as well as the percentage of the total frequency of all preposition-pronoun contractions occurring in the corpus is calculated. To determine whether the distribution of the unaccented postprepositional pronouns in the corpus may be considered linguistically significant and, in consequence, may establish a basis for the revision of the inflectional paradigms in (1)–(5), a number of quantitative procedures has been performed. In our ongoing study, the distribution of genitive and accusative feminine singular, as well as plural postprepositional pronouns has been analyzed. In future work, frequencies of the remaining pronouns will be examined, and, based on the evaluation results, the correctness of the inflectional paradigms in (1)–(5) will be confirmed or challenged.

- (1) Inflectional paradigm of the third person masculine human pronoun *on* 'he':

CASE	SINGULAR	PLURAL
nominative	<i>on</i> (—) / — (—)	<i>oni</i> (—) / — (—)
genitive	<i>jego</i> (go) / <i>niego</i> (-ń)	<i>ich</i> (—) / <i>nich</i> (—)
dative	<i>jemu</i> (mu) / <i>niemu</i> (—)	<i>im</i> (—) / <i>nim</i> (—)
accusative	<i>jego</i> (go) / <i>niego</i> (-ń)	<i>ich</i> (—) / <i>nich</i> (—)
instrumental	<i>nim</i> (—) / <i>nim</i> (—)	<i>nimi</i> (—) / <i>nimi</i> (—)
locative	— (—) / <i>nim</i> (—)	— (—) / <i>nich</i> (—)

- (2) Inflectional paradigm of the third person masculine animal pronoun *on* 'it':

CASE	SINGULAR	PLURAL
nominative	<i>on</i> (—) / — (—)	<i>one</i> (—) / — (—)
genitive	<i>jego</i> (go) / <i>niego</i> (-ń)	<i>ich</i> (—) / <i>nich</i> (—)
dative	<i>jemu</i> (mu) / <i>niemu</i> (—)	<i>im</i> (—) / <i>nim</i> (—)
accusative	<i>jego</i> (go) / <i>niego</i> (-ń)	<i>je</i> (—) / <i>nie</i> (—)
instrumental	<i>nim</i> (—) / <i>nim</i> (—)	<i>nimi</i> (—) / <i>nimi</i> (—)
locative	— (—) / <i>nim</i> (—)	— (—) / <i>nich</i> (—)

- (3) Inflectional paradigm of the third person masculine inanimate pronoun *on* 'it':

CASE	SINGULAR	PLURAL
nominative	<i>on</i> (—) / — (—)	<i>one</i> (—) / — (—)
genitive	<i>jego</i> (go) / <i>niego</i> (-ń)	<i>ich</i> (—) / <i>nich</i> (—)
dative	<i>jemu</i> (mu) / <i>niemu</i> (—)	<i>im</i> (—) / <i>nim</i> (—)
accusative	<i>jego</i> (go) / <i>niego</i> (-ń)	<i>je</i> (—) / <i>nie</i> (—)
instrumental	<i>nim</i> (—) / <i>nim</i> (—)	<i>nimi</i> (—) / <i>nimi</i> (—)
locative	— (—) / <i>nim</i> (—)	— (—) / <i>nich</i> (—)

- (4) Inflectional paradigm of the third person feminine pronoun *ona* 'she':

CASE	SINGULAR	PLURAL
nominative	<i>ona</i> (—) / — (—)	<i>one</i> (—) / — (—)
genitive	<i>jej</i> (—) / <i>niej</i> (—)	<i>ich</i> (—) / <i>nich</i> (—)
dative	<i>jej</i> (—) / <i>niej</i> (—)	<i>im</i> (—) / <i>nim</i> (—)
accusative	<i>ją</i> (—) / <i>nią</i> (—)	<i>je</i> (—) / <i>nie</i> (—)
instrumental	<i>nią</i> (—) / <i>nią</i> (—)	<i>nimi</i> (—) / <i>nimi</i> (—)
locative	— (—) / <i>niej</i> (—)	— (—) / <i>nich</i> (—)

- (5) Inflectional paradigm of the third person neuter pronoun *ono* 'it':

CASE	SINGULAR	PLURAL
nominative	<i>ono</i> (—) / — (—)	<i>one</i> (—) / — (—)
genitive	<i>jego</i> (go) / <i>niego</i> (—)	<i>ich</i> (—) / <i>nich</i> (—)
dative	<i>jemu</i> (mu) / <i>niemu</i> (—)	<i>im</i> (—) / <i>nim</i> (—)
accusative	<i>je</i> (—) / <i>nie</i> (—)	<i>je</i> (—) / <i>nie</i> (—)
instrumental	<i>nim</i> (—) / <i>nim</i> (—)	<i>nimi</i> (—) / <i>nimi</i> (—)
locative	— (—) / <i>nim</i> (—)	— (—) / <i>nich</i> (—)

	dlań 'for.PR'	doń 'to.PR'	nań 'on.PR'	weń 'in.PR'	zeń 'with.PR'	odeń 'from.PR'	przezeń 'by.PR'	poń 'after.PR'	zań 'behind.PR'	przedeń 'in front of.PR'	Total, %
[...]											[...] [...]
nom, n, sg											0 0.00 %
gen, n, sg	3	16			16	1					36 3.02 %
dat, n, sg											0 0.00 %
acc, n, sg			13	6			32		2		53 4.45 %
instr, n, sg											0 0.00 %
loc, n, sg											0 0.00 %
nom, n, pl											0 0.00 %
gen, n, pl		5									5 0.42 %
dat, n, pl											0 0.00 %
acc, n, pl				1			1				2 0.17 %
instr, n, pl											0 0.00 %
loc, n, pl				1							1 0.08 %
nom, f, sg											0 0.00 %
gen, f, sg	5	15			4	1					25 2.06 %
dat, f, sg											0 0.00 %
acc, f, sg			5	4			10				19 1.59 %
instr, f, sg											0 0.00 %
loc, f, sg											0 0.00 %
nom, f, pl											0 0.00 %
gen, f, pl	1	1			2	1					5 0.42 %
dat, f, pl											0 0.00 %
acc, f, pl			2				1				3 0.25 %
instr, f, pl											0 0.00 %
loc, f, pl			1								1 0.08 %
Total	101	219	377	101	93	23	250	1	27	1	1.193
%	8.47%	18.36%	31.60%	8.47%	7.80%	1.93%	20.96%	0.08%	2.26%	0.08%	100%

Table 1: The distribution of selected postprepositional unaccented pronouns in the IPI PAN Corpus